

Application-Specific Soft Processor Optimization through Instruction-Level Frequency Scaling

Sorin Dan Cotofana

Computer Engineering Laboratory
Faculty of Electrical Engineering, Mathematics and Computer Science
Delft University of Technology (TU Delft), The Netherlands

Overview

- Motivation
- FPGA Soft Processors as Accelerators
- Application-Specific Soft Processors
 - ISA Customization and Optimization
 - Instruction-Level Frequency Scaling
- Conclusions

Motivation

Estimated Chip Design Cost, by Process Node, Worldwide, 2011

Process Node	Design cost (\$M)	Mask cost (\$M)	Embedded software (\$M)	Yield ramp-up cost (\$M)
28/22-nm	~100	~10	~10	~10
32-nm	~80	~10	~10	~10
45-nm	~60	~10	~10	~10
65-nm	~45	~10	~10	~10
90-nm	~30	~10	~10	~10
130-nm	~15	~10	~10	~10
180-nm	~5	~10	~10	~10

TU Delft

3

Leveraging FPGA MPSoCs

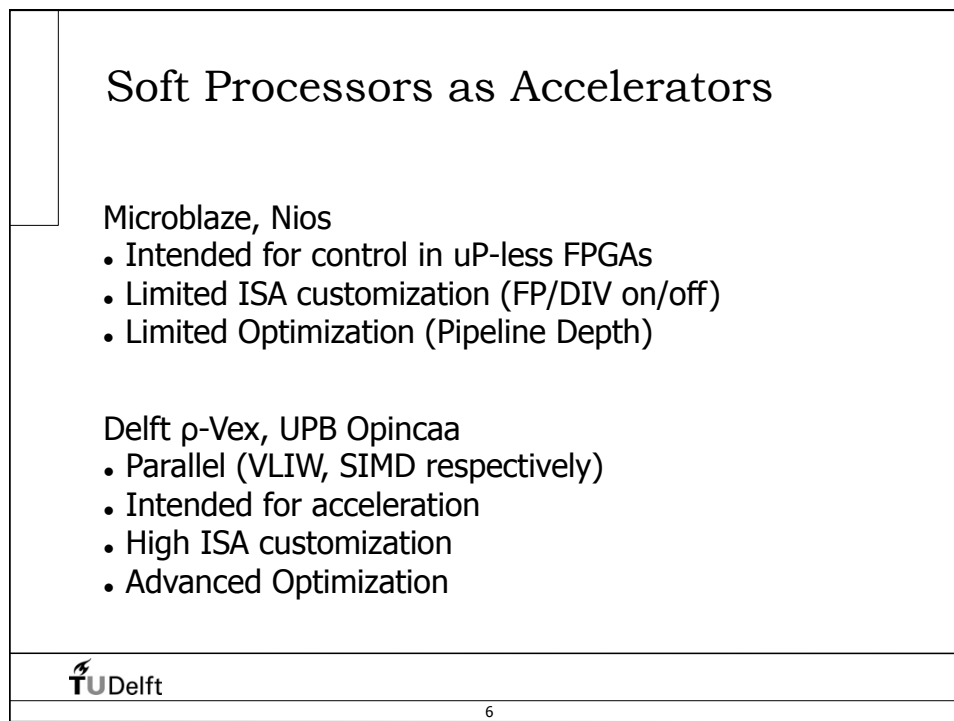
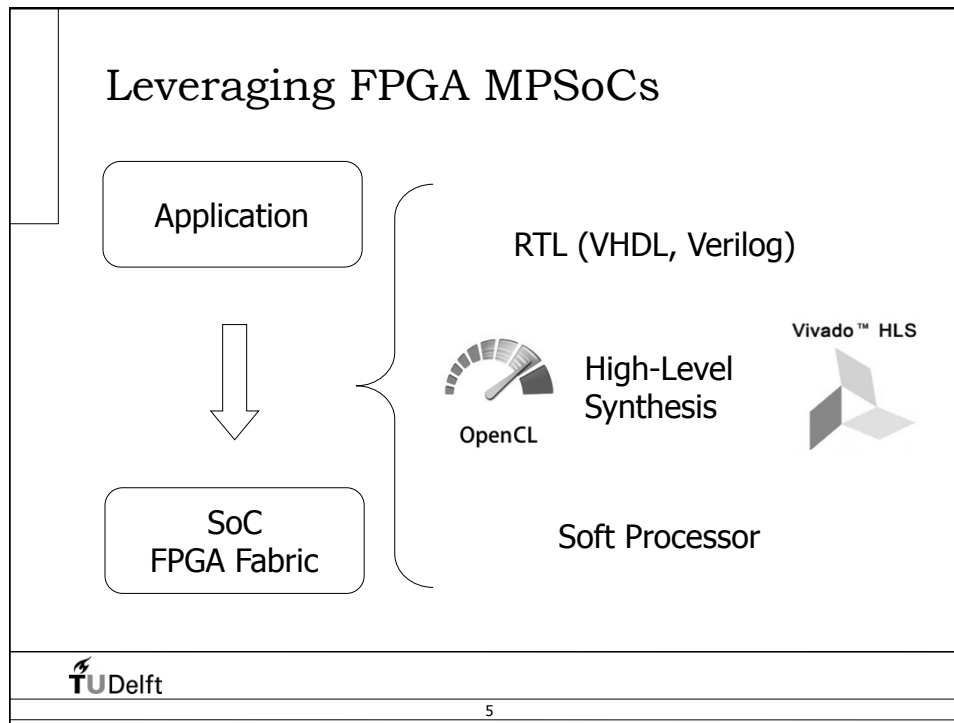
Application

Crunching Data Control Flow

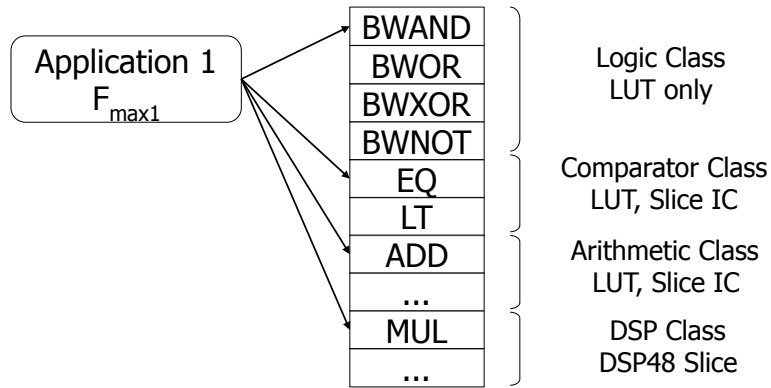
SoC FPGA Fabric e.g, ZYNQ™ SoC Processors

TU Delft

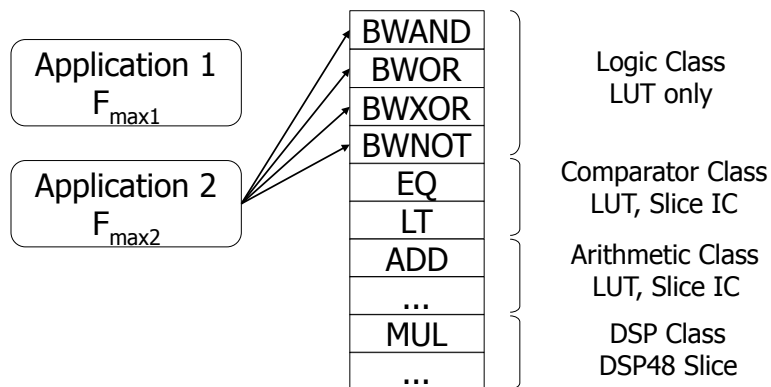
4

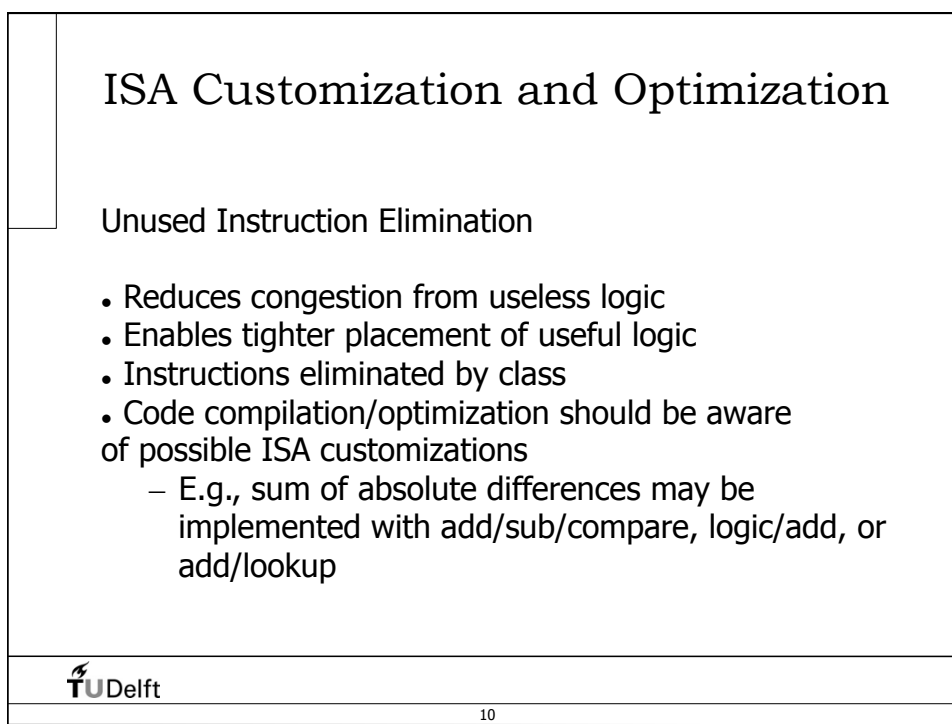
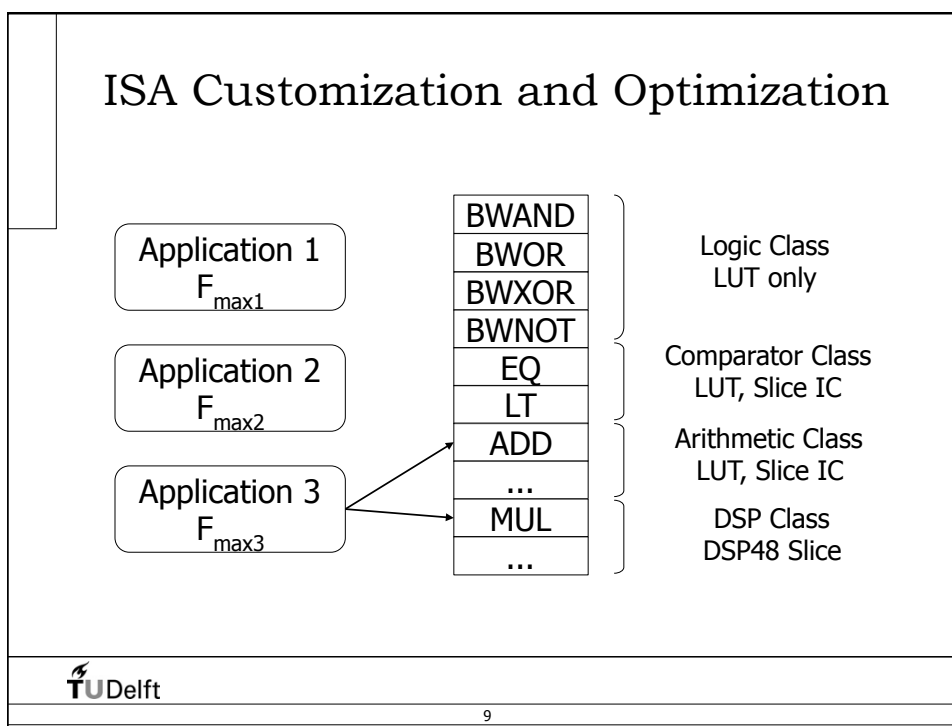


ISA Customization and Optimization



ISA Customization and Optimization





ISA Customization and Optimization

E.g., 128 SIMD ISA Customization on Zynq 7020

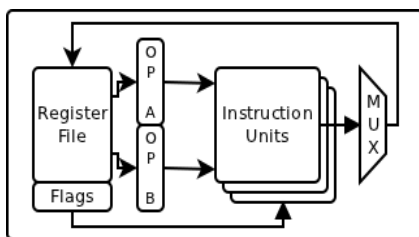
- Full ISA, $F_{\max} = \sim 80$ MHz (highly congested)
- Shift instruction removed, $F_{\max} = \sim 125$ MHz
- Add/Mul only, $F_{\max} = \sim 200$ MHz

SIMD Lane Width Customization

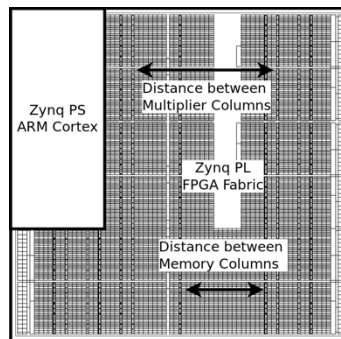
- Full ISA, SIMD lane width 128 (congestion)
- Add/Mul only, SIMD lane width 220 (DSP48)
- No Mul/Logic/Shift, SIMD lane width 276 (Block RAM)

Instruction-Level Frequency Scaling

- Instructions occur with varying rates
- Longest path delay dictates frequency
- No correlation between logic delay of instruction and occurrence rate



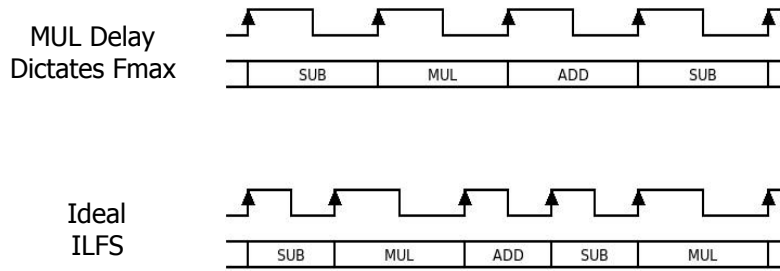
E.g. 128 SIMD on Zynq,
MUL incurs high routing delay



Instruction-Level Frequency Scaling

ILFS – decouple longest delay and top frequency

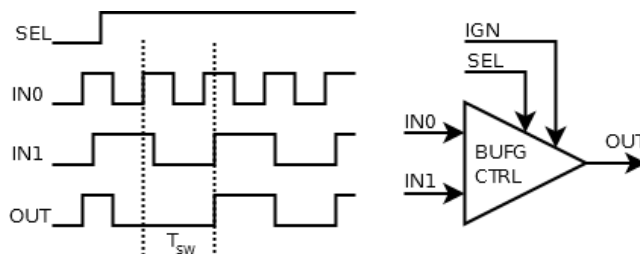
- E.g., 128 SIMD on Zynq 7020, F_{\max} 125 MHz
- MUL paths ignored, $F_{\max} \sim 160$ MHz



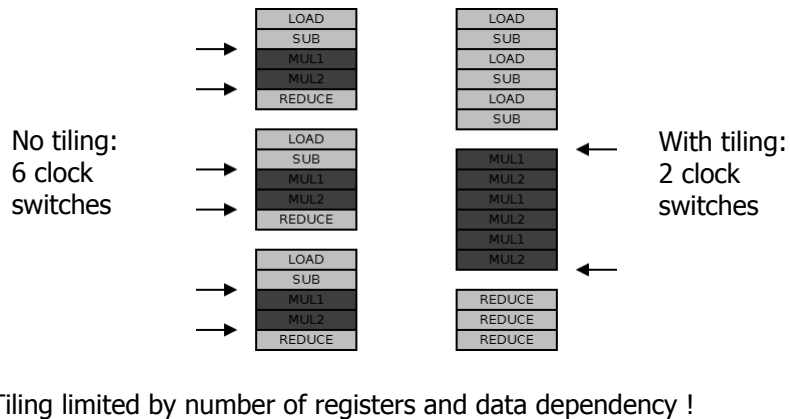
Instruction-Level Frequency Scaling

Zynq Implementation

- Uses BUFGCTRL clock multiplexer
- Incurs multiplexing latency at every clock frequency switch



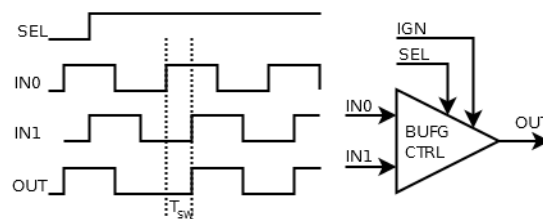
Instruction-Level Frequency Scaling



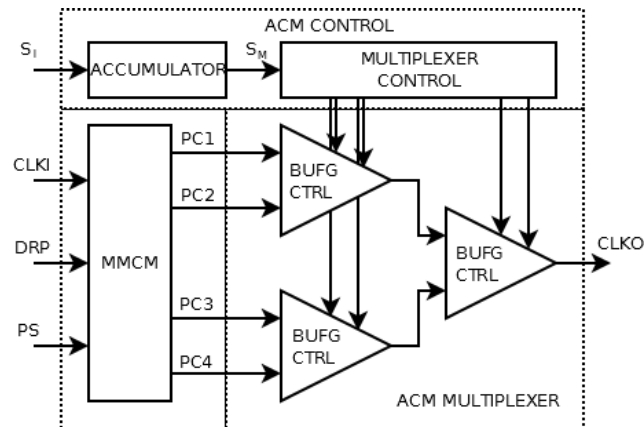
Instruction-Level Frequency Scaling

BUFGCTRL predictable when clocks have same F

- T_{SW} equals phase offset (source to destination)
- Phase offset known, T_{SW} is controlled stretch
- Continuous switch == continuous stretch
- Zero-latency T extension by T_{SW}



Instruction-Level Frequency Scaling



Instruction-Level Frequency Scaling

Automatic ILFS Configuration

- Profile application, determine instruction occurrence rates
- Apply separate timing constraints to paths pertaining to each instruction
- Tighten constraints for common instructions, loosen for rare instructions, synthesize
- Repeat while timing constraints met

Conclusion

- Soft processors an alternative to HLS or RTL
- ISA Customization enables application-specific soft processors, higher F_{\max} , wider SIMD
- Instruction-Level Frequency Scaling possible with existing FPGAs, enables higher F_{\max} even under difficult routing

Open Issues

- Customization-aware compiling
- ILFS per-instruction optimization difficult with current FPGA synthesis tools

Acknowledgement

Electronic Devices, Circuits and Architectures
Department, Politehnica University of Bucharest
(<http://www.dcae.pub.ro/en/>)

Opincaa Team:

- Radu Hobincu
- Călin Bîră
- Vlad Popescu
- Alexandru Gheolbănoiu
- Lucian Petrică (lucian.petrica@upb.ro)